

## REMARKS

### I. Introduction

In response to the Office Action dated April 24, 2003, no claims have been cancelled, amended or added. Claims 1-57 remain in the application. Re-examination and re-consideration of the application, as amended, is requested.

### II. Prior Art Rejections

#### A. The Office Action Rejections

In paragraphs (2)-(3) of the Office Action, claims 1, 20, and 39 were rejected under 35 U.S.C. §103(a) as being unpatentable over Bayer et al., U.S. Patent No. 5,202,987 (Bayer) in view of Tsuchida et al., U.S. Patent No. 6,026,394 (Tsuchida). In paragraph (4) of the Office Action, claims 1, 20, and 39 were rejected under 35 U.S.C. §103(a) as being unpatentable over Bhattacharya et al., U.S. Patent No. 5,797,000 (Bhattacharya). In paragraph (5) of the Office Action, claims 2-3, 21-22, and 40-41 were rejected under 35 U.S.C. §103(a) as being unpatentable over Bhattacharya as applied to claims 1, 20, and 39 in view of Hintz et al., U.S. Patent No. 5,222,235 (Hintz). In paragraph (6) of the Office Action, claims 4-6, 23-25, and 42-44 were rejected under 35 U.S.C. §103(a) as being unpatentable over Bhattacharya as applied to claims 1, 20, and 39 and in view of Bordonaro et al., U.S. Patent No. 5,307,485 (Bordonaro). In paragraphs (8)-(9) of the Office Action, claims 1, 20, and 39 were rejected under 35 U.S.C. §102(e) as being anticipated by Garth et al., U.S. Patent No. 6,272,486 B1 (Garth). However, in paragraph (10) of the Office Action, claims 12-19, 31-38, and 50-57 were indicated as being allowable if rewritten in independent form to include the base claim and any intervening claims.

Applicants' attorney acknowledges the indication of allowable claims, but respectfully traverse these rejections.

#### B. Applicants' Independent Claims

Applicants' independent claims 1, 16 and 30 are directed to loading data into a data store connected to a computer. Independent claim 1 is representative and comprises the steps of:

- identifying memory constraints;

- identifying processing capabilities; and

- determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities.

### C. The Bayer Reference

Bayer describes a high flow-rate synchronizer/scheduler apparatus for a mutiprocessor system during program run-time, which comprises a connection matrix for monitoring and detecting computational tasks which are allowed for execution containing a task map and a network of nodes for distributing to the processors information or computational tasks detected to be enabled by the connection matrix. The network of nodes possesses the capability of decomposing information on a pack of allocated computational tasks into messages of finer sub-packs to be sent toward the processors, as well as the capability of unifying packs of information on termination of computational tasks into a more comprehensive pack. A method of performing the synchronization/scheduling in a multiprocessor system of this apparatus is also described.

### D. The Tuschida Reference

Tsuchida describes a database management system for executing database operations in parallel by a plurality of nodes and a query processing method. The database management system contains a decision management node for deciding a distribution node for retrieving information so as to analyze a query received from an application program, generate a processing procedure for processing the query, and execute the process, and a join node for sorting, merging, and joining the information retrieved by the distribution node. When the query process is executed, the distribution node decided by the decision management node retrieves the information to be processed and the join node decided by the decision management node also obtains the result for the query from the retrieved information. The query result is outputted from an output node and transferred to the application program.

### E. The Bhattacharya Reference

Bhattacharya describes a method of performing a parallel join operation on a pair of relations R1 and R2 in a system containing P processors organized into Q clusters of P/Q processors each. The system contains disk storage for each cluster, shared by the processors of that cluster, together with a shared intermediate memory (SIM) accessible by all processors. The relations R1 and R2 to be joined are first sorted on the join column. The underlying domain of the join column is then partitioned into P ranges of equal size. Each range is further divided into M subranges of progressively decreasing size to create MP tasks T.sub.m,p, the subranges of a given

range being so sized relative to one another that the estimated completion time for task  $T_{sub.m,p}$  is a predetermined fraction that of task  $T_{sub.m-1,p}$ . Tasks  $T_{sub.m,p}$  with larger time estimates are assigned (and the corresponding tuples shipped) to the cluster to which processor  $p$  belongs, while tasks with smaller time estimates are assigned to the SIM, which is regarded as a universal cluster (cluster 0). The actual task-to-processor assignments are determined dynamically during the join phase in accordance with the dynamic longest processing time first (DLPT) algorithm. Each processor within a cluster picks its next task at any given decision point to be the one with the largest time estimate which is owned by that cluster or by cluster 0.

#### F. The Hintz Reference

Hintz describes a reorganization method of DB2 data files exploring parallel processing, and asynchronous I/O to a great extent. It includes means to estimate an optimum configuration of system resources, such as storage devices (DASD devices), memory, and CPUs, etc, during reorganizations. The method mainly consists of four components, (1) concurrent indexing, (2) concurrent unloading of data file partitions, (3) efficient reloading of DB2 data pages and DB2 space maps, and (4) means to reduce access constraints to the DB2 recovery table.

#### G. The Bordonaro Reference

Bordonaro describes a system and method for merging a plurality of sorted lists using multiple processors having access to a common memory in which  $N$  sorted lists which may exceed the capacity of the common memory are merged in a parallel environment. Sorted lists from a storage device are loaded into common memory and are divided into a number of tasks equal to the number of available processors. The records assigned to each task are separately sorted, and used to form a single sorted list. A multi-processing environment takes advantage of its organization during the creation of the tasks, as well as during the actual sorting of the tasks.

#### H. The Garth Reference

Garth describes a method, apparatus, and article of manufacture for a computer-implemented building indexes system. Indexes are built for a database that is stored in a data storage device coupled to a computer. An amount of available memory is determined. An amount of memory for use in transmitting data between extract, sort, and index build tasks is determined. Then, a number of sort tasks to be used to build indexes is determined based on the determined

amount of available memory, the determined amount of memory for use in transmitting data between tasks, and task memory requirements.

I. Applicants' Claimed Invention Is Patentable Over The References

Applicants' attorney respectfully submits that Applicants' claimed invention is patentable over the references. Specifically, Applicants' attorney asserts that the references do not teach or suggest the limitations recited in Applicants' independent claims 1, 20 and 39.

With regard to the rejections of claims 1, 20 and 39 under 35 U.S.C. §102(e) based on U.S. Patent No. 6,272,486 B1 (Garth), the present application has been amended to claim continuation-in-part status from U.S. Patent No. 6,272,486 B1 (Garth), thereby eliminating the patent as a prior art reference.

With regard to the rejections based on Bayer and Tsuchida, the Office Action states the following:

3. Claims 1, 20, and 39 are rejected under 35 U.S.C. 103(a) as being unpatentable over Bayer et al. (US Pat 5,202,987) in view of Tsuchida et al. (US Pat 6,026,394).

Regarding claims 1, 20, and 39, Bayer et al. disclose a method of loading data into a data store connected to a computer, the method comprising the steps of:

identifying memory constraints (col. 1, lines 13 - 15, memory and processors are operations bottleneck and col. 5, lines 52 - 56, memory is constrained or limited through factors such as physical shared storage, network access or processor distribution, and common memory space);

identifying processing capabilities (col. 1, lines 17 - 27, synchronization activities are controlled by algorithm, which depends on processing power and col. 5, lines 24 - 31 requires the number of processors and capabilities of each processor, which entail processing capabilities); and

determining a number of load (col. 14, lines 25 - 33, loading capacity being part of task map) to be started in parallel based on the identified memory constraints and processing capabilities (col. 7, lines 9 - 1).

Although Bayer et al. disclose the sort process being a mere tasks allocation to the processors (col. 1, lines 44 - 50), Tsuchida et al. have nevertheless further detailed the sort feature, which includes the step of determining a number of sort processes (col. 8, lines 50 - 51 disclose the fact that the sorting process depends on the number of node for join process. Col. 7, lines 54 - 57 show that the number of join nodes for performing merge process can be determined. Hence, number of sort processes is a known quantity).

Therefore, it is considered obvious to one of ordinary skill in the art, at the time the invention was made, to combine the sorting feature shown by Tsuchida et al. to the invention of Bayer et al. so that sort processing time, which is a factor in load balancing processes can be determined as part of system characteristics and

optimization purposes (col. 7, lines 58 - col. 8, line 35). Note that the sort steps shown by Tsuchida et al. are also parallel processes as claimed in the application (fig. 3, parallel pipeline operation).

Applicants' attorney disagrees. The cited portions of these references do not teach or suggest the limitation "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities."

For example, the cited portions are set forth below:

Bayer: Col. 1, lines 13-27

The coordination of multiple operations in shared memory multiprocessors often constitutes a substantial performance bottleneck. Process synchronization and scheduling are generally performed by software, and managed via shared memory. Execution of parallel programs on a shared-memory, speedup-oriented multiprocessor necessitates a means for synchronizing the activities of the individual processors. This necessity arises due to precedence constraints within algorithms: When one computation is dependent upon the result of other computations, it must not commence before they finish. In the general case, such constraints are projected onto an algorithm's parallel decomposition, and reflected as precedence relations among its execution threads.

Bayer: Col. 5, lines 24-31 (actually, 24-44)

In addition to a task map, the synchronizer/scheduler is supplied with the system configuration data. This includes such details as the number of processors, the capabilities of each processor (if processors are not a-priori identical), etc.

Given a set of enabled tasks, as well as processor availability data, the synchronizer/scheduler then performs scheduling of those tasks. Any non-random scheduling policy must rely upon some heuristics: Even when task execution times are known in advance, finding an optimal schedule for a program represented as a dependency graph is an NP-complete problem. Most scheduling heuristics are based on the critical path method, and thereby belong to the class of list scheduling policies; i.e., policies that rely on a list of fixed task priorities. List scheduling can be supported by the inventive scheme described herein, by embedding task priorities in the task map load-module submitted to the synchronizer/scheduler. Whenever an allocation takes place, the allocated tasks are those which have highest priorities amongst the current selection of enabled tasks.

Bayer: Col. 14, lines 25-33

Characterizing Parameters

The parameters characterizing a specific synchronizer/scheduler can now be summarized:

Loading Capacity:

The maximal size of a task map which can be loaded. This parameter is expressed in terms of quantity of tasks, and/or in terms of quantity of dependency connections.

Bayer: Col. 7, lines 9-11

The multiprocessor architecture is illustrated in FIG. 1. As can be seen, the parallel operation coordination subsystem (synchronizer/scheduler 10) forms an appendage to a conventional configuration of a shared-memory 12 and processors 14.

Tsuchida: Col. 8, lines 50-51

The slot sorting process is set so as to be executed by the nodes for join process.

Tsuchida: Col. 7, lines 54-57

Next, the method for deciding the number of join nodes for performing the N-way merge process will be explained with reference to FIG. 7. FIG. 7 is a schematic view for explaining the decision method for the number of join nodes.

Tsuchida: Col. 7, line 58 - col. 8, line 35

In FIG. 7, it is assumed that the data retrieval/distribution process is executed in the nodes # 1 to # 8 and the processing time in each node is the one shown at each of the numbers 300 to 305. In this example, the processing time 304 in the node # 5 is the maximum processing time. The slot sorting processing time can be driven from the number of nodes for join process N, predetermined system characteristics (CPU performance, disk unit performance, etc.), and database operation method. The performance characteristic (processing time  $E_s$ ) of the slot sorting process can be obtained generally from the following expression.

$$E_s = a/N + b*N + c \quad (1)$$

The N-way merge processing time ( $E_m$ ) and join processing time ( $E_j$ ) also can be obtained from the following expressions.

$$E_m = d/N + e*N + f \quad (2)$$

$$E_j = g/N + h*N + i \quad (3)$$

where, symbols a, d, and g indicate constants which are decided from system characteristics such as the number of rows, the number of pages, each operation unit time, and output time. Symbols b, e, and h are constants which are decided from system characteristics such as the communication time, and c, f, and i are constants which are decided from the other system characteristics.

According to this embodiment, to maximize the effect of the pipeline process, the number of nodes for join process is obtained as the number of assigned join nodes 350 so that the performance characteristic  $E_s$  of the slot sorting process becomes equal to the maximum processing time 304. When the number of assigned join nodes 350 is determined, the N-way merge processing time 320 and join processing time 330 can be estimated from the equations (2) and (3). The total of these processing times is the total processing time for a query. By deciding the number of join nodes in this manner and merging the data distributed in the data retrieval/distribution process successively and processing them simultaneously, the total processing time (response time from querying to output) can be shortened.

The descriptions set forth above do not teach or suggest the limitation "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities."

Instead, Bayer is directed to the scheduling of synchronized tasks, but tasks are assigned to processors based on availability, wherein a processor is allocated a new task immediately after it terminates the previous one. See, e.g., col. 15, lines 1-25. Moreover, the "Loading Capacity" referred to above as a "Characterizing Parameter," relates to the maximal size of the task map, wherein the task map is a data structure that identifies dependencies between tasks being performed, and is used as indicated to assign tasks to processors based on availability.

In addition, Tsuchida is directed a parallel processing database management system, wherein the configuration of the nodes, i.e., processes, is fixed. See, e.g., col. 2, line 59 - col. 3, line 49. Moreover, a decision means determines which (already started) nodes are to be used to perform the query in order to minimize the expected processing time.

As a result, neither of the Bayer or Tsuchida references teach or suggest "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities." Consequently, it cannot be said that the combination of Bayer and Tsuchida teaches or suggests, or renders obvious, the Applicant's independent claims.

With regard to the rejections based on Bhattacharya, the Office Action states the following:

4. Claims 1, 20, and 39 are rejected under 35 U.S.C. 103(a) as being unpatentable over Bhattacharya et al. (US Pat 5,797,000).

As per claims 1, 20, and 39, Bhattacharya et al. disclose a method of loading data into a data store connected to a computer, the method comprising the steps of: identifying memory constraints (col. 9, lines 6 - 7, memory becomes a constraint as its capacity is a contributing factor and is limited);

identifying processing capabilities (fig. 1, number of processors p, col. 4, line 41- col. 5, line 8, each processor is assigned with a specific number of tasks, hence indicating each limited capability); and

determining a number of load (col. 3, lines 1 - 3, join column domain and tuples are the load, which obviously must be known in order for them to be partitioned and transferred among the cluster, col. 3, lines 7 - 18), and sort processes (col. 2, line 62 - col. 3, line 6 disclose various method of parallel sort process in which merge join is one example. Since the actual tasks assigned to the processors are determined during the join phase, which is part of the sort process as shown above, number of sort processes are hence inherently determined as well) to be started in parallel based on the identified memory constraints and processing capabilities (col. 1, lines 24 - 28, parallel tasks based on processing capabilities: col. 4, line 64 - col. 5, line 8, parallel sort processing: col. 2, lines 66 - col. 3, line 6).

Applicants' attorney disagrees. The cited portions of the reference do not teach or suggest the limitations "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities."

For example, the cited portions are set forth below:

Bhattacharya: Col. 9, lines 6-7 (actually 6-17)

Each processor 104 of the system 100 is allotted an equal portion 1/P of the memory capacity of universal cluster 108. In the initial portion of the transfer phase, for each processor p (104) of the system 100, the tasks T.sub.m,p corresponding to that processor and residing on a particular cluster 102 are transferred from that cluster to the universal cluster 108, beginning with the task T.sub.M,p having the smallest estimated completion time and progressing in order of increasing task size (i.e., decreasing m), until the allotted portion 1/P is filled (step 532). The remaining tasks T.sub.m,p for each processor p (104) are transferred to the cluster 102 owning the processor, unless they are already resident there (step 534).

Bhattacharya: Col. 4, line 41 - col. 5, line 8

Referring to FIG. 1, a multiprocessor system 100 incorporating the present invention includes P processors 104 organized into Q equal-size clusters 102, each cluster containing P/Q processors. Each processor 104 may be either a uniprocessor or a complex of tightly coupled processors (not separately shown) that, for the purposes of task assignment, are regarded as a single processor. Each cluster 102 also includes one or more direct access storage devices (DASD) 106, which are magnetic disk drives in the system 100 shown. Each processor 104 within a cluster 102 can access any storage device 106 in the same cluster, but cannot access any storage device in any other cluster 102. Processors 104 are interconnected to one another as well as to a single intermediate memory (SIM) 108, to which each processor has access. SIM 108 is also referred to herein as the universal cluster, or cluster 0. In addition to the memory 108 and storage devices 106 shown, each processor 104 also has its own main memory (not separately shown). In the case of a processor 104 comprising a tightly coupled processor complex, such main memory would be shared by the processors of the complex. The elements shown in FIG. 1 are conventional in the art, as are the interconnections between these elements.

Processors 104 are used for the concurrent parallel execution of tasks making up database queries, as described below. A query may originate either from one of the processors 104 or from a separate front-end query processor as described in the concurrently filed application of T. Borden et al., Ser. No. 08/148,091, now U.S. Pat. No. 5,495,606. As further described in that application, within each cluster 102 the query splitting and scheduling steps described below may be performed by an additional processor or processors (not shown) similar to processors 104; such additional processors would not be counted among the P/Q processors 104 per complex 102 to which tasks are assigned.

Bhattacharya: Col. 3, lines 1-3 (actually col. 2, line 66 - col. 3, line 6)

In a parallel sort merge join, the relations to be joined are first sorted, in parallel, within their clusters 102 (FIG. 1). In a naive parallel sort merge join, the



underlying join column domain might be partitioned into  $P$  ranges of equal size, and the tuples transferred accordingly among the clusters 102. However, given a nonuniform distribution of tuples across the underlying domain, there is no guarantee that the amount of join phase work will be equal.

Bhattacharya: Col. 3, lines 7-18 (actually lines 7-24)

In accordance with the present invention, each of the  $P$  ranges is further divided into a relatively small number  $M$  of components, creating  $MP$  tasks  $T_{sub.m,p}$  in all. These components intentionally have nonequal task time estimates. For example, a reasonable approach would be to partition the tasks so that the estimated completion time of a task  $T_{sub.m,p}$  is half that of the previous task  $T_{sub.m-1,p}$ . Assuming that the quadratic output term dominates the task time estimates, this can be done by partitioning the tasks in such a manner that the extent of the range of a given task  $T_{sub.m,p}$  (to which the number of tuples in the task is roughly proportional) is  $1/\sqrt{2}$  times the number of tuples in task  $T_{sub.m-1,p}$ . FIGS. 10A and 10B show an example of such a partitioning. FIG. 10A shows estimated task times as a function of  $m$  and  $p$ , and FIG. 10B shows actual task times, also as a function of  $m$  and  $p$ . The latter may be different from the former, and will not be known until the join phase, when the tasks are actually performed.

Bhattacharya: Col. 1, lines 24-28

This invention relates generally to a method of performing a parallel query in a multiprocessor environment and, more particularly, to a method for performing such a query with load balancing in an environment with shared disk clusters, shared intermediate memory or both.

The descriptions set forth above do not teach or suggest the limitation "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities."

Instead, Bhattacharya is directed to a parallel join operation, wherein tasks are assigned to processors based on the partitioning of the domain of the join column into  $P$  ranges and the partitioning of each range into  $M$  subranges to create  $MP$  tasks. In this context,  $P$  is based on the number of processors and  $M$  is based on the underlying domain, e.g., each subrange has an extent in the underlying column domain that is  $1/\sqrt{2}$  that of the preceding subrange, and the subranges of a given range being so sized relative to one another, such that the estimated completion time for task  $T_{mp}$  is a predetermined fraction that of task  $T_{m-1,p}$ . In Bhattacharya, tasks  $T_{mp}$  with larger time estimates are assigned to the cluster to which processor  $P$  belongs, while tasks with smaller time estimates are assigned to a universal cluster (e.g., cluster 0). The actual task-to-processor assignments are determined dynamically during the join phase in accordance with the dynamic longest processing time first (DLPT) algorithm. Each processor within a cluster picks its next task

at any given decision point to be the one with the largest time estimate which is owned by that cluster or by cluster 0.

As a result, Bhattacharya does not teach or suggest the limitations "determining a number of load and sort processes to be started in parallel based on the identified memory constraints and processing capabilities." Consequently, it cannot be said that Bhattacharya renders the Applicants' independent claims obvious.

Hintz and Bordonaro fail to overcome the deficiencies of Bhattacharya. Recall that Hintz was cited only against dependent claims 2-3, 21-22 and 40-41, while Bordonaro was cited only against dependent claims 4-6, 23-25 and 42-44. Moreover, Hintz was cited only for determining a number of build processes based on the number of sort processes, and for teaching that the number of sort processes does not exceed a number of indexes to be built, while Bordonaro was cited only for teaching that the number of load processes does not exceed a number of partitions to be loaded, and that the load and sort processes directly dependent on memory constraints. None of these teachings are relevant to the limitations of Applicants' independent claims.

Thus, Applicants submit that independent claims 1, 20 and 39 are allowable over the references. Further, dependent claims 2-19, 21-38 and 40-57 are submitted to be allowable over the references in the same manner, because they are dependent on independent claims 1 and 12, respectively, and thus contain all the limitations of independent claims 1 and 12. In addition, dependent claims 4-9, 11-25 and 27-44 recite additional novel elements not shown by the references.

### III. Conclusion

In view of the above, it is submitted that this application is now in good order for allowance and such allowance is respectfully solicited.

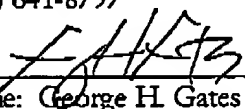
Should the Examiner believe minor matters still remain that can be resolved in a telephone interview, the Examiner is urged to call Applicants' undersigned attorney.

Respectfully submitted,

GATES & COOPER LLP  
Attorneys for Applicants

Howard Hughes Center  
6701 Center Drive West, Suite 1050  
Los Angeles, California 90045  
(310) 641-8797

Date: July 24, 2003

By:   
Name: George H. Gates  
Reg. No.: 33,500

GHG/

G&C 30571.279-US-01